# Experiences with a Multi-Protocol network monitor

Andrew Moore

Intel Research Fellow

Computer Laboratory

University of Cambridge

In collaboration with:

Ian Pratt, Jon Crowcroft, James Hall, Tim Granger, Derek McAuley, Dina Papagiannaki, among others.

# Contents

- λ  one-slider on Projects @ Cambridge
- λ  Nprobe/GRIDprobe Monitor
- λ  Experiences with...
  - λ  Behaviour example
  - λ  Content example
  - λ  Visualizer example
- λ  Where next...

# Projects at Cambridge on Monitoring

- λ **Nprobe/GRIDprobe/Xenoprobe**

    Computer Laboratory

    (Collectively known as *probe or (star)probe)

- λ **CoMo**

    Intel Research Cambridge + friends

    Planning for convergence, also taking in Hyperion (U.Mass)

# GRIDprobe Objective

- λ (Nprobe prototyped/grew-into GRIDprobe)
- λ Scalable Monitoring Architecture (tool building)
  - λ 1 & 10Gbps and viable strategy for 40Gbps
- λ Multi-protocol monitoring
  - λ Understand network and application behaviour
    At the same time.
- λ Originally University of Cambridge & Marconi (RIP)
- λ Now Cambridge with association from Intel
- λ Duration Oct. 2002 – Sep. 2005

# Status

- λ Several working test deployments (1Gbps)
- λ Prototype for 10 Gbps
- λ Code base is planned for a public release
- λ Experience with the dataset/database/dataware-house issues
- λ Adding new protocol modules
- λ Using Experience to drive next architecture

# Where does this tool fit in?

"When you have a hammer, every problem looks like a nail."

- λ We want current network data
  - λ High-resolution timer
  - λ High-speed (current deployment: 1 Gbps)
- λ We want to collect enough information to see the interaction between layers
- λ We want to use commodity (no custom) hardware to maximize deployment and minimize cost

# Nprobe: our current implementation

λ Current Nprobe system performs full line-rate capture on commodity hardware

λ Nprobe is a multi-protocol monitor: collecting network, transport & application data

λ Nprobe processes network, transport & application layers  to provide compression as well as extracting useful information (e.g., application features)

# What we are NOT

λ We are not just some IDS – they do a few things that superficially look the same – ultimately these things are not the same.

λ We want to collect as much as possible – they want to collect the minimum and to compare as quick as possible.

λ we want to interpret the full application – they want to string-match then move on.

# What is the problem?

In a perfect world:

Cheaply (using commodity PCs)

Record 1, 10, (MAXINT) Gbps

Full duplex

Onto disk

With minimal loss

Ouch!

Not as bad as all that: its not a
perfect world

# How do we do it?

"Discard is the most effective compression."

**Be selective** (for an http example)

1. Remove redundant header information
2. Temporally compress header information
3. Extract http transactions from data stream
4. Remove (or summarise) uninteresting information (consider the use)

# How else?

"A problem shared is a problem halved."

Split the workload

λ    among CPUs

λ    among machines

Problems?

λ    Complex filter
     design made
     easier using
     ongoing
     measurements

Limitations?

# Limitations Abilities

λ  Host-host data must be less than or equal to the capacity of a single monitor (CPU)

λ  No monitoring DataTAG 10 Gbps host-host experiments

λ  For ISP and dial-up or cable modem last-miles, as well as with (UK) academics with 100 Mbps to the desktop, this approach works

λ  Target deployment has nx10,000 of flows and the monitor is close to the server or close to the client (on access/choke-points).

# Example 1: Modelling TCP Connections

λ Dynamic model of TCP connection activity

    λ Input from probe-collected data

        λ Packet timings

        λ Packet header data

        λ Higher level protocol activity

    λ Output — identified, differentiated and quantified

        λ Network times

        λ TCP Artefacts

        λ Application delays

# Causative Associations

λ Probe sees TCP packets traveling to host and those returning

    λ Arriving packets

        λ Modify host TCP state

        λ Cause work to be done

        λ Trigger transmissions — causative associations

        λ Drive model

    λ Departing packets

        λ Verify/modify model

        λ Are arrivals at peer

# Example 1:

λ ACK packet arrives during slow start

  ∖ sender's congestion window expands

  ∖ releases flight of data segment packet(s)

λ Data segment N arrives {N mod 2 = 0}

  ∖ ACK released

λ Data segment N transmitted

  ∖ data segment N+1 released

λ HTTP request arrives

  ∖ First packet of response released (after delay)

# *Partial* Round Trip Times

λ Probe can be anywhere

    λ Hence deal in *p*RTTs)

    λ Can glue them together

# Lags, Delays and *p*RTTs

λ **For causative associations**
- λ lag = *p*RTT + delay
- λ If no delay:
  - λ *p*RTT = lag
- λ If delay:
  - λ interpolate *p*RTT
  - λ Calculate delay
- λ Model informs

# *p*RTT Drawbacks/Restrictions

λ  Only works in slow-start, thus relies on longer data flows

λ  relies on implementation "inside knowledge" fortunately only a few implementations (BSD derivative, Linux derivative, Microsoft derivative)

# Results — Live Traffic

λ  All HTTP traffic to BBC news server from University site

  λ  24 Hour trace

  λ  Results for period 1130 – 1350

    λ  Expect load increase as users browse during lunch break

λ  Independent of local load

λ  Look at SYN re-transmissions

# Local Load — Live Traffic

# Server Delays and *p*RTTs

# Server *p*RTTs — Live Traffic

# Probabiliy-of-Retransmission

# Server Delay and   p-SYN Retransmission

# Behaviour Summary

λ By observing a combination of TCP and HTTP protocols simultaneously, we:

  λ determine the type of load-shedding this website uses.

  λ understand diminished performance in the face of no local network effects.

  λ Draw conclusions on the impact this approach (to load-shedding) has on persistent vs. non-persistent (compared with a nominal 25%, this site had less than 5% persistent)

# GRIDprobe Visualization Tools

λ Reads the stored format – the stored format is already partially processed

λ Extracts features of interest (timing, packet headers,...)

λ Constructs relationship trees for (web) pages

λ provides:

λ interactive data plotter (ala gnuplot++)

λ tcp connection plotter

λ web transaction plotter

# Why?

λ **Aids understanding**

    λ of observed behaviour,

    λ and trends

λ **Teaching tool**

λ **Debugging tool**

λ **Maintains relationship between layers**

    λ tcp/ip ... http/html ... coarse statistics

## Analysis log

```
13 User_Agents seen
31 Pages seen (5 > 5 objects) 11 with delays (3 > 5 objects) 0 persistent
24 clients
8 servers
152 connections 129 objects, 0 on persistent connections 22 URLs
24 no request connections - 9 incorporated
14 objects - 28 conns 2 servers not in trees
Total connections: 199
Rejected:
Rejected by tfilter: 2
SYNs not seen  2
Total transactions 0
Valid 0
Rejected:
User Agents and inferred TCP implementations:
Servers and inferred TCP implementations:
Log saved to /home/jhall/nprobe/logs/A.rep.0112.r200.NOBT5_treelog


Ranks

 Client-by-#objects
 Client-by-#bytes
 Client-by-#conns
 Client-by-#pages
 Server-by-#objects
 Server-by-#bytes
 Server-by-#conns
 Server-by-#pages
 Server-by-pdt-per-object-(>50b/P)
 Server-by-pdt-per-object-(>50b/P)-less-delays
 Server-by-pdt-per-object-(>50b/P)-85%
 Server-by-pdt-per-object-(>50b/P)-85%-less-delays
 Pages-(>5ob)-by-frequency
 Pages-(>5ob)-by-download-time-per-object
 Pages-(>5ob)-by-download-time-per-object-less-delays
```

```
rtmr, 0 fin-delay
rtmr, 3 fin-delay


jects) 0 persistent



s 22 URLs
```

**Quit**  **By_Id**  By_Type  **Shrink**  **Clear**  |  **Draw**

Delays/object over time (1)

Delays/connection over time (1)

No. attempts/connection (1)

Retransmits/connection over time (1)

All Page Downloads (3)
　　All page downloads (31) (1)
　　All page downloads - pages <= 5 (26) (1)
　　All page downloads - pages > 5 (5) (1)

　DUP SYN (20)
　HIGH ACK (19)
0/1 Uconns incorporated (1)
0/4 Uconns incorporated (1)
1/2 Uconns incorporated (1)
1/3 Uconns incorporated (1)
7/14 Uconns incorporated (1)
C Retransmit (1)
Client dup syns - clients (7)
Client dup syns - conns (24)
Client dup syns - no req - clients (1)
Client dup syns - no req - conns (1)
Client dup syns - not connected - clients (2)
Client dup syns - not connected - conns (5)
Client early loss - clients (1)
Client early loss - conns (1)
Client rtmts (1)
Client transaction invalid/error (1)
DUP FIN (3)
Dummy Transaction (2)
Fail registered (1)
Large simpage first resp (4)
Link not found in referring object (3)
MULTIACK (18)

Menu (Draw):

Connections
Page
All Pages

**Page Times**
**85% Page Times**
**Page Times less delays**
**85% Page Times less delays**

**Page Times per-object**
**85% Page Times per-object**
**Page Times less delays per-object**
**85% Page Times less delays per-object**

Retransmits
Retransmits-P-over-1s
Retransmits-P-over-5s
Retransmits-P-over-10s
Retransmits-P-over-60s

Delays
**Page-delays**
**Page-delays-per-object**
Attempts

**Cancel**

Right window text (partially visible):

mr, 0 fin-delay
r, 3 fin-delay

ts) 0 persistent

22 URLs

oads - pages <= 5 (26)'

np_plot /home/jhall/nprobe/logs/A.rep.0112.r200.NOBT5 - All page times

Quit  Style  sets  TS  HIST  PDF  CDF  Smooth  Print  Data  Add

Page
downloads

0

15000

12500

10000

7500

5000

2500

2.5          5.0          7.5          10.0          12.5          15.0

elapsed time s

All page times

Quit    Prev    Next    Detail    Conns    Grab    Redo    Reload

**WebClient xx.xx.xx.xx**
Ref Tree root
http://www.bbc.co.uk:80/weather/ukweat
r/east/index.shtml conn 151287

Start 11:59:11.707324 17/5/2001
Duration 1802.029ms

22 Connections
(max. 5/4/4 concurrent/eff. concurrent
22 transactions
(0 invalid 0 failed)

objects linked: 22  unlinked: 0

Object types:-

image/gif 0/19
image/jpeg 0/2
text/html 0/1
type unknown 22/0

Link types:-

------ Claimed followed link 0
········ Claimed inline link 2
········ In-line link 19

Server return codes:-
200  OK 22

User agents:-
(a) Mozilla/4.77 [en] (X11; U; Linux
2.2.19-6.2.1smp i686) 22

0    50    100    150    200

>1 http://www.bbc.co.uk:80/weather/

1 (<--1) (a) GET http://www.bbc.co.uk:80/weather/ukweather/ea
(b) C (200)

0 (151287) PO: 0/0/0    212.58.224.36
1.0/1.1    0
5

+182

1 (151292) PO: 1/1/1    212.58.224.36
1.0/1.1    23

Quit  Prev  Next  Detail  Conns  Grab  Redo  Reload

WebClient xx.xx.xx.xx
Ref Tree root
http://www.bbc.co.uk:80/weather/ukweat
r/east/index.shtml conn 151287

Start 11:59:11.707324 17/5/2001
Duration 1802.029ms

22 Connections
(max. 5/4/4 concurrent/eff. concurrent
22 transactions
 (0 invalid 0 failed)

objects linked: 22   unlinked: 0

Object types:-

███████  image/gif 0/19
███████  image/jpeg 0/2
███████  text/html 0/1
███████  type unknown 22/0

Link types:-

----- Claimed followed link 0
...... Claimed inline link 2
......  In-line link 19

Server return codes:-
 200  OK 22

User agents:-
 (a) Mozilla/4.77 [en] (X11; U; Linux
    2.2.19-6.2.1smp i686) 22

450    500    550    600    650    700    750

(+526) C: 553/18 S: 35661/32
er/images/gui/bg_ffffcc.gif                        +491
(b) C (200) 131/1
(+209) C: 5/2/6 S: 378/5
+180
5 (<-1) (a) GET http://www.bbc.co.uk:80/weather/images/gui/dotsv_ffffcc.gif
212.58.224.36                                      (b) C (200) 49/1
1.0/1.1    478                                      (+195) C: 575/6 S: 295/5
490                    +158
4 (<-1) (a) GET http://www.bbc.co.uk:80/weather/images/gui/dotsh_ffffcc.gif
212.58.224.36                                      (b) C (200) 75/1
1.0/1.1    479                                      (+185) C: 575/6 S: 321/5
489                    +162
3 (<-1) (a) GET http://www.bbc.co.uk:80/weather/maps/easttemp012.jpg
212.58.224.36                                      (b) C (200) 9580/7
1.0/1.1    480                                      (+231) C: 568/8 S: 9
488                    +166
5 (<-1) (a) GET http://www.bbc.co.uk:80/weather/images/gui/button
299) PO: 3/3/3        212.58.224.36               (b) C (200) 152/1
1.0/1.1    528                                      (+198) C: 575/6 S: 399/5
548                    +124
7 (<-1) (a) GET http://www.bbc.
6 (151300)  PO: 3/1/1   212.58.224.36
1.0/1.1    666
676  8 (<-1) (a) GET http://www.bb
7 (151301)  PO: 3/1/1   212.58.224.36
1.0/1.1    675
681  9 (<-1) (a) GET http://w
8 (151303)  PO: 3/2/2   212.58.224.36
1.0/1.1    687
701  10 (<-1) (a) GET ht
9 (151305)  PO: 3/3/3   212.58.224.36
1.0/1.1    712
718

10 (151306)  PO: 3/2/2   212
1.0.

11 (151307)  PO: 3/2/2   2
1

12 (151308)  PO: 3/3/3

13 (151309)  P

**Reference Tree for WebClient xx.xx.xx.xx**

Quit   Prev   Next   Time   Conns   Grab   Redo   Reload

**WebClient xx.xx.xx.xx**

Ref Tree root
http://www.bbc.co.uk:80/weather/ukweat
r/east/index.shtml conn 151287

Start 11:59:11.707324 17/5/2001
Duration 1802.029ms

22 Connections
(max. 5/4/4 concurrent/eff. concurrent)
22 transactions
 (0 invalid 0 failed)

objects linked: 22  unlinked: 0

Object types:-

image/gif 0/19
image/jpeg 0/2
text/html 0/1
type unknown 22/0

Link types:-

Claimed followed link 0
Claimed inline link 2
In-line link 19

Server return codes:-
 200  OK 22

User agents:-
 (a) Mozilla/4.77 [en] (X11; U; Linux
    2.2.19-6.2.1smp i686) 22

0 4        188 188      254 308     430 430     452 452     456 456     463 464     477 478
4 4        188 193        308 308     430 430     452 454     456 458     464 464     478 479
4 5          193 232       308 309     430 430     454 454     458 463     464 476     479 480
5 9          232 237       309 309     430 438     454 454     463 463     476 476     480 48
  9 187        237 237     309 310     438 442     454 455     463 463     476 476     482 4
 187 187        237 250     310 333     442 452     455 456     463 463     476 477     483
 187 186        250 254     333 430     452 452     455 456     463 463     477 477     48

**>1** http://www.bbc.co.uk:80/weather/

**1 (<--1)** (a) GET http://www.bbc.co.uk:80/weather/ukweather/east/index.shtml
         (b) C (200) 35531/28

212.58.224.36
.0/1.1    0

5 +182              **2 (<-1)** (a) GET http://www.bbc.co.uk:80/weather/images/gui/bg_ffffc
                         (b) C (200) 131/1
2) PO: 1/1/1   212.58.224.36              (+209) C: 572/6 S: 378/5
               1.0/1.1    232
                   250        +16080

                                   **2 (151296)  PO: 1/1/1**   212.58.224.36
                                                               1.0/1.1    478

                                   **3 (151297)  PO: 2/2/1**   212.58.224.36
                                                               1.0/1.1    479

                                   **4 (151298)  PO: 3/3/1**   212.58.224.36
                                                               1.0/1.1    480

                                                           **5 (151**

Quit  Prev  Next  Detail  Conns  Grab  Redo  Reload

00    450    500    550    600    650    700    750

**WebClient xx.xx.xx.xx**
Ref Tree root
http://www.bbc.co.uk:80/weather/ukweat
r/east/index.shtml conn 151287

Start 11:59:11.707324 17/5/2001
Duration 1802.029ms

22 Connections
(max. 5/4/4 concurrent/eff. concurrent
22 transactions
 (0 invalid 0 failed)

objects linked: 22  unlinked: 0

Object types:-

image/gif 0/19
image/jpeg 0/2
text/html 0/1
type unknown 22/0

Link types:-
Claimed followed link 0
Claimed inline link 2
In-line link 19

Server return codes:-
 200  OK 22

User agents:-
 (a) Mozilla/4.77 [en] (X11; U; Linux
 2.2.19-6.2.1smp i686) 22

(+526) C: 553/18 S: 35661/32
weather/images/gui/bg_ffffcc.gif
(b) C (200) 131/1
(+209) C: 5/2/6 S: 378/5
+180
**5 (<-1)** (a) GET http://www.bbc.co.uk:80/weather/images/gui/dotsv_ffffcc.gif
1/1/1   212.58.224.36                                          (b) C (200) 49/1
1.0/1.1   478                                                 (+195) C: 575/6 S: 295/5
490
**4 (<-1)** (a) GET http://www.bbc.co.uk:80/weather/images/gui/dotsh_ffffcc.gif
2/2/1   212.58.224.36                                    +158  (b) C (200) 75/1
1.0/1.1   479                                                 (+195) C: 575/6 S: 321/5
489
**3 (<-1)** (a) GET http://www.bbc.co.uk:80/weather/maps/easttemp012.jpg
3/3/1   212.58.224.36                              +162  (b) C (200) 9580/7
1.0/1.1   480                                            (+231) C: 568/8
488                                                       +146
**5 (<-1)** (a) GET http://www.bbc.co.uk:80/weather/images/gui/bu
(151299)  PO: 3/3/3   212.58.224.36                 (b) C (200) 152/1
1.0/1.1   528                                        (+198) C: 575/6 S: 399/
548                                                      +124
**7 (<-1)** (a) GET http://www
6 (151300)  PO: 3/1/1   212.58.224.36
1.0/1.1   666    676
**8 (<-1)** (a) GET http://w
7 (151301)  PO: 3/1/1   212.58.224.36
1.0/1.1   675    681
**9 (<-1)** (a) GET htt
8 (151303)  PO: 3/2/2   212.58.224.36
1.0/1.1   687    701
**10 (<-1)** (a) GE
9 (151305)  PO: 3/3/3   212.58.224.36
1.0/1.1   712  718

10 (151306)  PO: 3/2/2

11 (151307)  PO: 3/2/

12 (151308)  PO: 3/

13 (151309

Reference Tree for WebClient xx.xx.xx.xx

Packet list for connection #151296

James pretty picture of connection #151296

Quit    Show    Server-Imp    Cient-Imp    Print    Grab

Re-draw    Sim

```
 0  478.569 >    0 +   1 w:32120 S    trig=none (0) lag=none
 1  482.684 <    0 +   1 ack1 w:24616 SA   trig=SYNACK (40)
 2  483.474 >    1 +   0 ack1 w:32120 A    trig=ACK (20) lag=
 3  490.134 >    1 + 575 ack1 w:32120 AP   trig=first (40000)
 4  494.339 <    1 +   0 ack576 w:24041 A   trig=del ACK (10
 5  649.000 <    1 + 295 ack576 w:24616 AP   trig=first (1000
 6  649.002 <  296 +   1 ack576 w:24616 AF   trig=none (0) la
 7  650.567 >  576 +   0 ack296 w:31856 A    trig=del ACK (10
 8  650.569 >  576 +   0 ack297 w:31856 A    trig=FINACK (80
 9  669.713 >  576 +   1 ack297 w:31856 AF   trig=none (0) la
10  673.791 <  297 +   0 ack577 w:24616 A    trig=FINACK (8
```

#151296 Client xx.xx.xx.xx  Server 212.58.224.36 port 80

 Non-persistent

User agent Mozilla/4.77 [en] (X11; U; Linux 2.2.19-6.2.1smp i686)
Server Apache/1.3.14 (Unix)

Start 11:59:12.185893 17/5/2001 duration 195.000ms  C: 575/6 mss 1460 S:
295/5 mss 1460

Server imp Imp=Gen IWF=0 SSTGT=1   Client imp Imp=Gen IWF=0 SSTGT=1

server
sequence

client
sequence

575
550
525
500
475
450
425
400
375
350
325
300
275
250
225
200
175
150
125
100
75
50
25
0.0

575
550
525
500
475
450
425
400
375
350
325
300
275
250
225
200
175
150
125
100
75
50
25
0

500    525    550    575    600    625    650    ms

pkts in
flight

server
prtt ms

client
prtt m

10 (151306) PO: 3/2/2
1.0/1.1

11 (151307) PO: 3/2/2
1.0/1

12 (151308) PO: 3/3/3
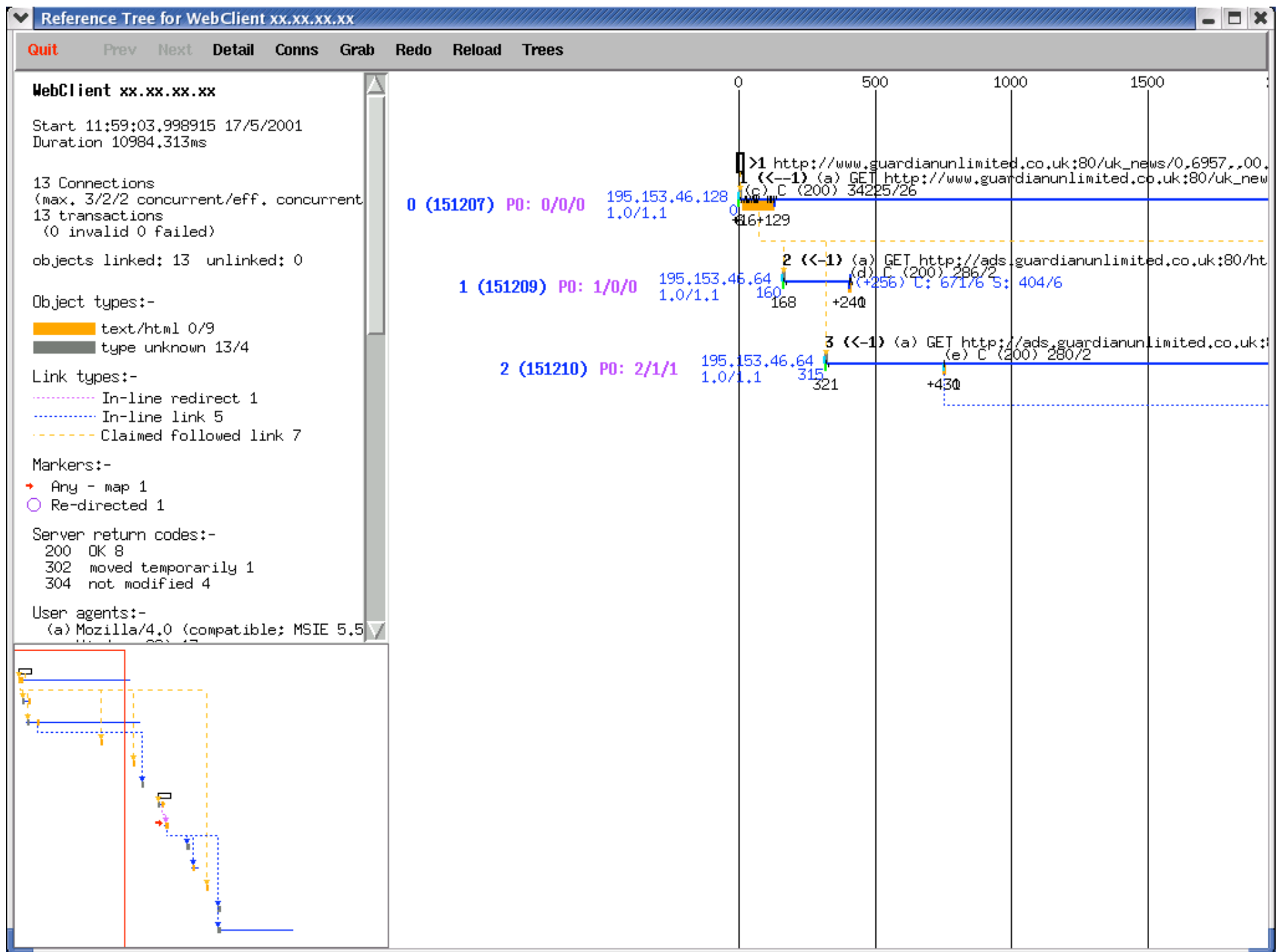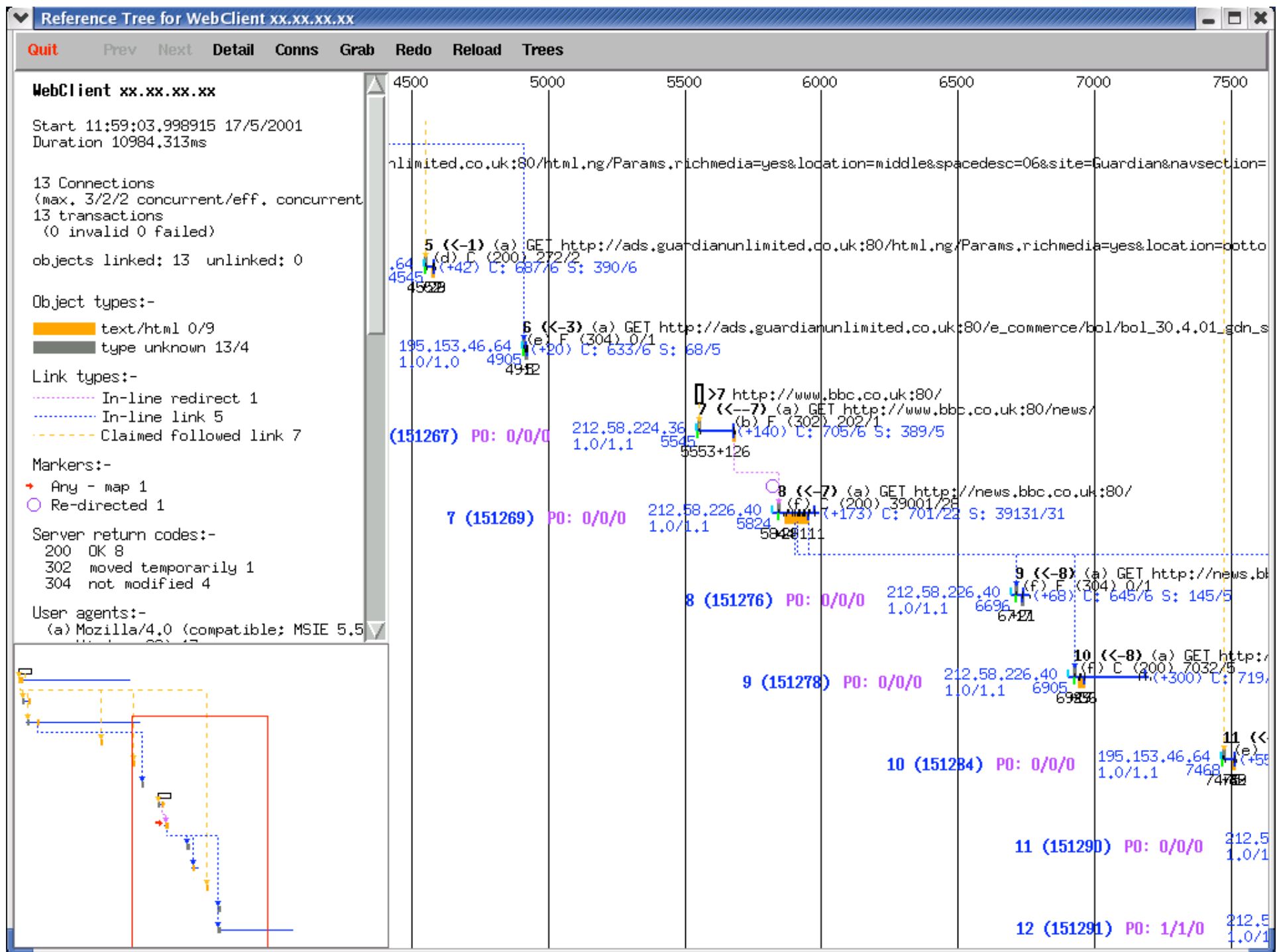1.0

13 (151309) PO: 3/

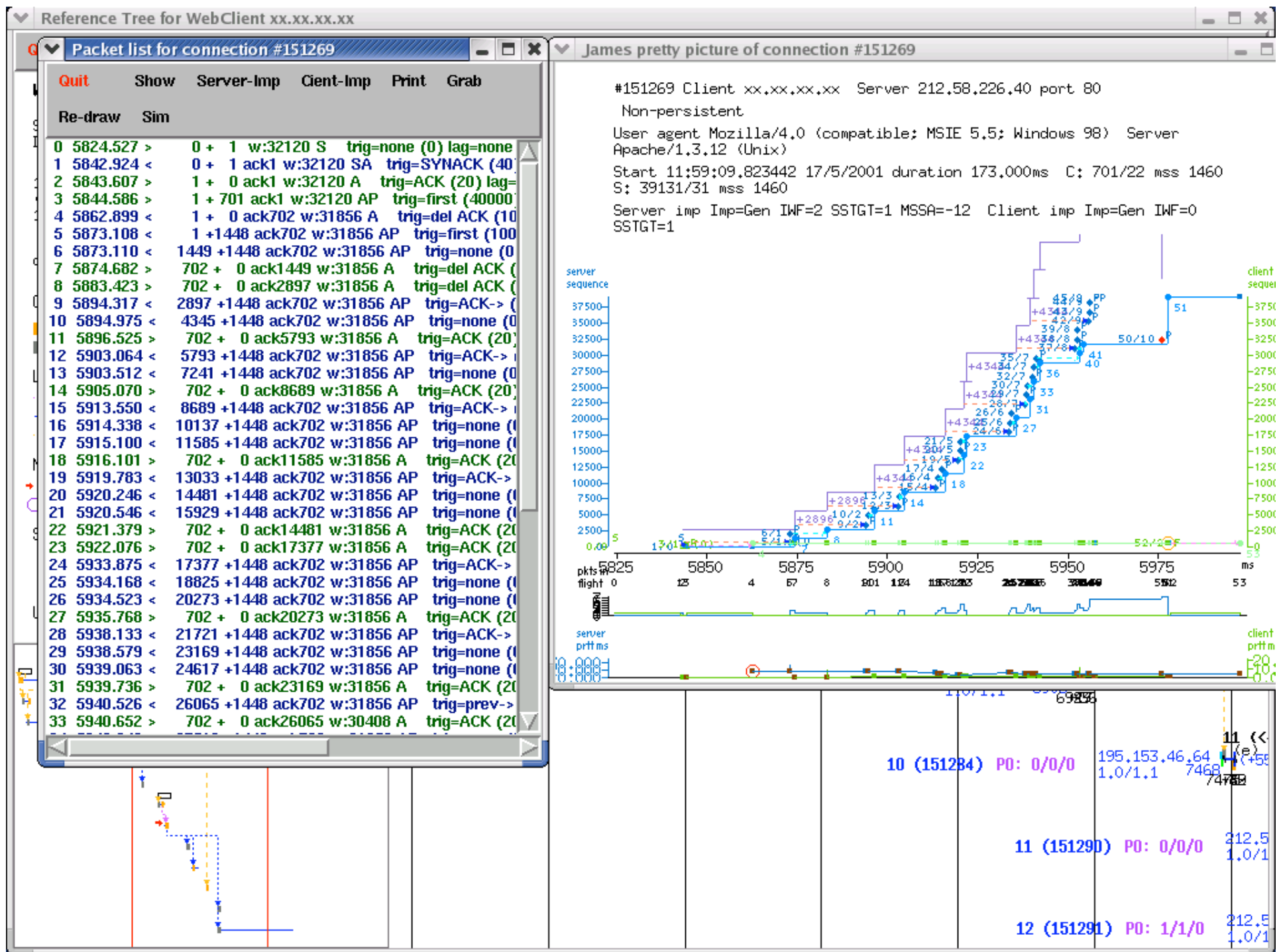**Quit**  **By_Id**  By_Type  **Shrink**  **Clear**  **Draw**

No request - conns (2)
No server data packets (3)
No server data packets - conns (25)
No server response - clients (1)
No server response - conns (1)
No transactions (9)
Not connected - clients (4)
Not connected - conns (22)
Poss fast-rtmt (3)
Redirection (3)
  xx.xx.xx.xx
  xx.xx.xx.xx
  xx.xx.xx.xx
S MSS adjusted (5)
S Retransmit (5)
S short segment (13)
SEG AFTER RELEASE (2)
SEG END EXCEEDS HIGH RELEASE (1)
Seen (24)
Server delayed fin - clients (2)
Server delayed fin - conns (3)
Server early loss - clients (1)
Server early loss - conns (1)
Server late loss - clients (2)
Server late loss - conns (4)
Server long rtmt - clients (1)
Server long rtmt - conns (1)
Server rtmts (7)
Set_IW - rel_highseq = 87802742 snd_nxt = 87793982
Set_IW - too many initial releases (1)
Unlinked object (10)
Zero client IWF (111)
Zero server IWF (90)
Zero window (32)
rtmt already ACKd (4)

End

Connections

**Page**

**All Pages**

Page Times

85% Page Times

Page Times less delays

85% Page Times less delays

Page Times per-object

85% Page Times per-object

Page Times less delays per-object

85% Page Times less delays per-object

Retransmits

Retransmits-P-over-1s

Retransmits-P-over-5s

Retransmits-P-over-10s

Retransmits-P-over-60s

Delays

Page-delays

Page-delays-per-object

Attempts

**Cancel**
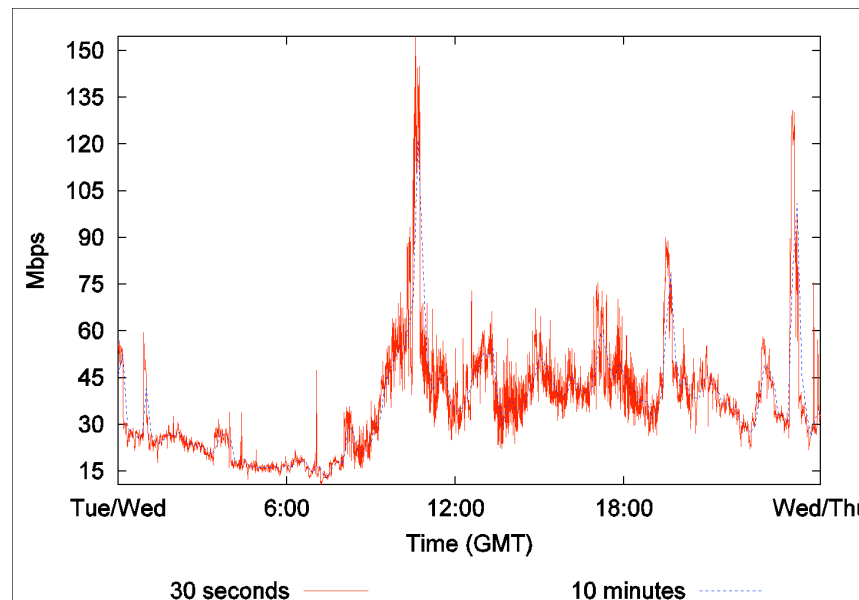
1335000/images/_1335019

1335000/images/_1335755

~~~~~~

~~~~~~

**Quit**   Prev   Next   **Detail**   **Conns**   **Grab**   **Redo**   **Reload**   **Trees**

**WebClient xx.xx.xx.xx**

Start 11:59:03.998915 17/5/2001
Duration 10984.313ms

13 Connections
(max. 3/2/2 concurrent/eff. concurrent)
13 transactions
 (0 invalid 0 failed)

objects linked: 13   unlinked: 0

Object types:-

▮ text/html 0/9
▮ type unknown 13/4

Link types:-

........ In-line redirect 1
........ In-line link 5
------- Claimed followed link 7

Markers:-
→ Any - map 1
○ Re-directed 1

Server return codes:-
 200  OK 8
 302  moved temporarily 1
 304  not modified 4

User agents:-
 (a) Mozilla/4.0 (compatible; MSIE 5.5

0                500              1000             1500

>1 http://www.guardianunlimited.co.uk:80/uk_news/0,6957,,00.
1 (<--1) (a) GET http://www.guardianunlimited.co.uk:80/uk_new
(c) C (200) 34225/26

**0 (151207) PO: 0/0/0**    195.153.46.128
1.0/1.1    0
16+129

2 (<--1) (a) GET http://ads.guardianunlimited.co.uk:80/ht
(d) C (200) 286/2
(+256) C: 6/1/6 S: 404/6

**1 (151209) PO: 1/0/0**    195.153.46.64
1.0/1.1    160
168        +240

3 (<--1) (a) GET http://ads.guardianunlimited.co.uk:
(e) C (200) 280/2

**2 (151210) PO: 2/1/1**    195.153.46.64
1.0/1.1    315
321        +430

Quit   Prev   Next   Detail   Conns   Grab   Redo   Reload   Trees

**WebClient xx.xx.xx.xx**

Start 11:59:03.998915 17/5/2001
Duration 10984.313ms

13 Connections
(max. 3/2/2 concurrent/eff. concurrent)
13 transactions
 (0 invalid 0 failed)

objects linked: 13  unlinked: 0

Object types:-

▬ text/html 0/9
▬ type unknown 13/4

Link types:-

......... In-line redirect 1
......... In-line link 5
--------- Claimed followed link 7

Markers:-
↦ Any - map 1
○ Re-directed 1

Server return codes:-
 200  OK 8
 302  moved temporarily 1
 304  not modified 4

User agents:-
 (a) Mozilla/4.0 (compatible; MSIE 5.5

4500     5000     5500     6000     6500     7000     7500

nlimited.co.uk:80/html.ng/Params.richmedia=yes&location=middle&spacedesc=06&site=Guardian&navsection=

**5 (<-1)** (a) GET http://ads.guardianunlimited.co.uk:80/html.ng/Params.richmedia=yes&location=botto
.64 (d) F (200) 272/2
4545 (+42) C: 6877/6 S: 390/6
4522

**6 (<-3)** (a) GET http://ads.guardianunlimited.co.uk:80/e_commerce/bol/bol_30.4.01_gdn_s
195.153.46.64 (e) F (304) 0/1
1.0/1.0   4905 (+20) C: 633/6 S: 68/5
4912

>7 http://www.bbc.co.uk:80/
**7 (<--7)** (a) GET http://www.bbc.co.uk:80/news/
(b) F (302) 202/1
**(151267)  PO: 0/0/0**   212.58.224.36 (+140) C: 705/6 S: 389/5
1.0/1.1   5545
5553+126

○ **8 (<-7)** (a) GET http://news.bbc.co.uk:80/
(f) F (200) 39001/28
**7 (151269)  PO: 0/0/0**   212.58.226.40 (+173) C: 701/22 S: 39131/31
1.0/1.1   5824
5848111

**9 (<-8)** (a) GET http://news.bb
(f) F (304) 0/1
**8 (151276)  PO: 0/0/0**   212.58.226.40 (+68) C: 645/6 S: 145/5
1.0/1.1   6696
6721

**10 (<-8)** (a) GET http:/
(f) C (200) 7032/5
**9 (151278)  PO: 0/0/0**   212.58.226.40 (+300) C: 719.
1.0/1.1   6905
6956

**11 (<**
(e) (+55
**10 (151284)  PO: 0/0/0**   195.153.46.64
1.0/1.1   7468
7752

**11 (151290)  PO: 0/0/0**   212.5
1.0/1

**12 (151291)  PO: 1/1/0**   212.5
1.0/1

**Packet list for connection #151269**

Quit    Show    Server-Imp    Cient-Imp    Print    Grab

Re-draw    Sim

```
 0 5824.527 >      0 +  1 w:32120 S    trig=none (0) lag=none
 1 5842.924 <      0 +  1 ack1 w:32120 SA   trig=SYNACK (40
 2 5843.607 >      1 +  0 ack1 w:32120 A    trig=ACK (20) lag=
 3 5844.586 >      1 + 701 ack1 w:32120 AP   trig=first (40000
 4 5862.899 <      1 +  0 ack702 w:31856 A    trig=del ACK (10
 5 5873.108 <      1 +1448 ack702 w:31856 AP   trig=first (100
 6 5873.110 <   1449 +1448 ack702 w:31856 AP   trig=none (0
 7 5874.682 >    702 +  0 ack1449 w:31856 A   trig=del ACK (
 8 5883.423 >    702 +  0 ack2897 w:31856 A   trig=del ACK (
 9 5894.317 <   2897 +1448 ack702 w:31856 AP   trig=ACK-> (
10 5894.975 <   4345 +1448 ack702 w:31856 AP   trig=none (0
11 5896.525 >    702 +  0 ack5793 w:31856 A   trig=ACK (20)
12 5903.064 <   5793 +1448 ack702 w:31856 AP   trig=ACK->
13 5903.512 <   7241 +1448 ack702 w:31856 AP   trig=none (0
14 5905.070 >    702 +  0 ack8689 w:31856 A   trig=ACK (20)
15 5913.550 <   8689 +1448 ack702 w:31856 AP   trig=ACK->
16 5914.338 <  10137 +1448 ack702 w:31856 AP   trig=none (0
17 5915.100 <  11585 +1448 ack702 w:31856 AP   trig=none (0
18 5916.101 >    702 +  0 ack11585 w:31856 A   trig=ACK (20
19 5919.783 <  13033 +1448 ack702 w:31856 AP   trig=ACK->
20 5920.246 <  14481 +1448 ack702 w:31856 AP   trig=none (0
21 5920.546 <  15929 +1448 ack702 w:31856 AP   trig=none (0
22 5921.379 >    702 +  0 ack14481 w:31856 A   trig=ACK (20
23 5922.076 >    702 +  0 ack17377 w:31856 A   trig=ACK (20
24 5933.875 <  17377 +1448 ack702 w:31856 AP   trig=ACK->
25 5934.168 <  18825 +1448 ack702 w:31856 AP   trig=none (0
26 5934.523 <  20273 +1448 ack702 w:31856 AP   trig=none (0
27 5935.768 >    702 +  0 ack20273 w:31856 A   trig=ACK (20
28 5938.133 <  21721 +1448 ack702 w:31856 AP   trig=ACK->
29 5938.579 <  23169 +1448 ack702 w:31856 AP   trig=none (0
30 5939.063 <  24617 +1448 ack702 w:31856 AP   trig=none (0
31 5939.736 >    702 +  0 ack23169 w:31856 A   trig=ACK (20
32 5940.526 <  26065 +1448 ack702 w:31856 AP   trig=prev->
33 5940.652 >    702 +  0 ack26065 w:30408 A   trig=ACK (20
```

**James pretty picture of connection #151269**

#151269 Client xx.xx.xx.xx  Server 212.58.226.40 port 80

Non-persistent

User agent Mozilla/4.0 (compatible; MSIE 5.5; Windows 98)  Server
Apache/1.3.12 (Unix)

Start 11:59:09.823442 17/5/2001 duration 173.000ms  C: 701/22 mss 1460
S: 39131/31 mss 1460

Server imp Imp=Gen IWF=2 SSTGT=1 MSSA=-12  Client imp Imp=Gen IWF=0
SSTGT=1

server sequence

client sequence

pkts #
flight

server prtt ms

client prtt ms

10 (151284)  P0: 0/0/0    195.153.46.64
                          1.0/1.1    7468

11 (151290)  P0: 0/0/0    212.5
                          1.0/1

12 (151291)  P0: 1/1/0    212.5
                          1.0/1

# Example 2: Experiences at a GRID Access Point

λ Research community of 1,000 on their own campus

λ Significant (unique) data provided by this site to the world community

λ One of three sites where data is continuously updated – so data is continuously transferred

between partners and downloaded by collaborators

# Traffic to/from access point

Total Link traffic



Each Direction

# Contrasting port and content based classification

| | Port-based | | Content-based | |
|---|---|---|---|---|
| | Packets | Bytes | Packets | Bytes |
| BULK | 49.97 | 45.00 | 65.06 | 64.94 |
| DATABASE | 0.03 | 0.03 | 0.84 | 0.96 |
| GRID | 0.03 | 0.07 | 0.00 | 0.00 |
| INTERACTIVE | 1.19 | 0.43 | 0.75 | 0.39 |
| MAIL | 3.37 | 3.62 | 3.37 | 3.62 |
| SERVICES | 0.07 | 0.02 | 0.29 | 0.28 |
| WWW | 19.98 | 20.40 | 26.50 | 27.60 |
| UNKNOWN | 28.36 | 30.44 | - | - |
| OTHER | - | - | 3.20 | 3.11 |

# Overheads vs. Accuracy

From Moore, Papagiannaki submission to IMC

(measures in packets)

| Method | UNKNOWN | Correctly Identified |
|---|---|---|
| Port | 29% | 71% |
| 1KB Signature | 24% | 74% |
| 1KB Protocol Parse | 19% | 81% |
| Important flow Assembly/Parse | 1% | 98% |
| Full Assemble/Parse | 0% | 100% |

# Classification Surprises

λ  Significant asymmetry to/from site

λ  Port-based classification was **so** wrong

λ  Considered the most important node for it's work yet,

λ  No GRID-application traffic!

It was all GRID web services or FTP traffic

(For the GRID community this was surprising, for the rest of us – less so.)

# Conclusions

λ   Our approach is sound: the implementation works and has been perfectly satisfactory for the environments into which we have deployed

λ   Clear avenues of development are available to us

λ   Further work with deployments will provide input to this work and provide data for other projects too

# What have we learnt?

λ Always things to improve:

   λ hardware

   λ optimization


λ Important to remember:

   λ compression of between 1:12 and 1:50 is achieved

   λ output data is all ready for off-line processing

# Next...

λ  Continue work to identify and quantify GRID and Internet applications:

e.g., Access GRID

λ  Evaluating 10 Gbps scheme

   λ  using test environments

   λ  considering deployment (10Gbps surprisingly uncommon)

# Enabling...

$\lambda$  Using of 10Gbps for the new UKLIGHT testbed – following the growth from implementation to deployment into full use.

$\lambda$  Current work enables us to assess suitability of Peer-2-peer algorithms for distributing data currently shared using FTP...

# Backup slides

# Architecture

λ Data driven model

λ Single-thread model to maximize efficiency

λ Avoid memory copy when practical

# Monitor / Hardware

Objective: use commodity hardware even the NIC

λ   modify the firmware rather than building hardware

λ   add/use time-stamping – currently 1_s

λ   perform filtering  on card with minimal overhead

λ   Current 1Gbps cards supported:

    λ   Alteon / 3Com 3c985B

    λ   SysKonnect sk98xx

# Monitor / Hardware Filtering

λ **Using hash of XOR of SRC/DEST as a selection criterion:**

- λ our approach works best when both directions of traffic are kept together.

λ **Work in progress**

- λ how often do we need to update filters?
- λ what can we optimize filters for:
  - λ filter size?
  - λ packet distribution?
  - λ equalizing flows? packets? bytes?
  - λ This problem is common to the load-sharing community

# Receive FIFO Implementation

Receive Buffer is not a FIFO

- This means that data-processing mechanisms can return data buffers when they are finished

- Out-of-order return allows easier handling of packet loss and packet reordering

- Discards packets when memory runs low

- Implemented to hang on to packets in case of potential use or reuse

# Monitor / Processing

Compress/discard where we can:

- λ  network, TCP and application layers can each have considerable temporal redundancy

- λ  application-specific reductions such as removing the data object from http transactions – we keep a fingerprint of the object so as we can recognise the same object even with different URLs

- λ  loss-less compression (lz, gzip, etc.)

# Nprobe modules

λ **Current**

    λ TCP and UDP on IP

    λ HTTP and HTML

    λ DNS

    λ FTP

λ **Past (deprecated)**

    λ NFS (v2)